**National Aeronautics and
Space Administration**

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

# Analytics Center Framework for
# Estimating the Circulation and Climate of the Ocean

**Thomas Huang**

thomas.huang@jpl.nasa.gov

Group Supervisor - Computer Science for Data-Intensive Applications
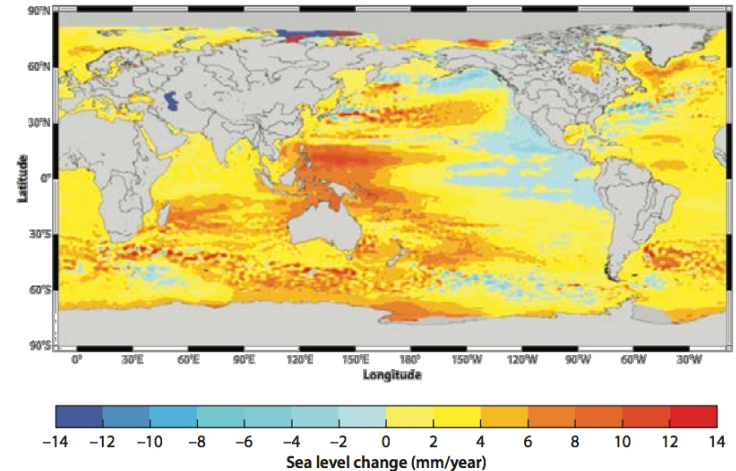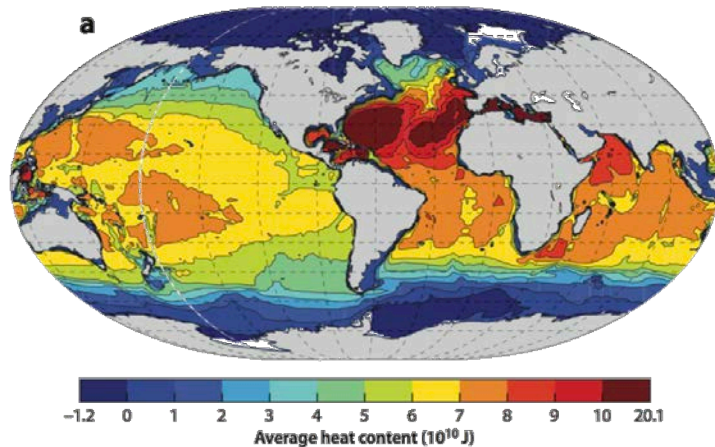
Strategic Lead - Interactive Data Analytics

Jet Propulsion Laboratory

California Institute of Technology

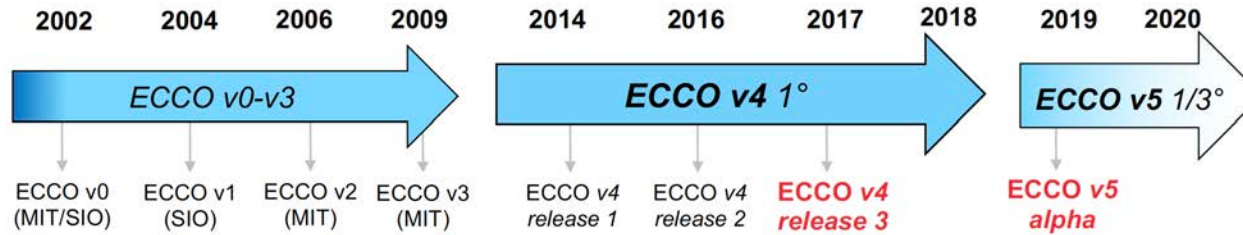4800 Oak Grove Drive, Pasadena, CA 91109-8099, U.S.A.

# What is ECCO?

- **Estimating the Circulation and Climate of the Ocean** (ECCO) is a consortium endeavors to produce the best possible estimates of ocean circulation and its role in climate
- Combining state-of-the-art ocean circulation models with global ocean and sea-ice data in a physically and statistically consistent manner
- ECCO products are being used in studies on ocean variability, biological cycles, coastal physics, water cycle, ocean-cryosphere interactions, and geodesy



a

−1.2   0   1   2   3   4   5   6   7   8   9   10   20.1
**Average heat content ($10^{10}$ J)**



−14   −12   −10   −8   −6   −4   −2   0   2   4   6   8   10   12   14
**Sea level change (mm/year)**

# ECCO Central Production Timeline

- ECCO v4 is the latest release of ocean state estimate
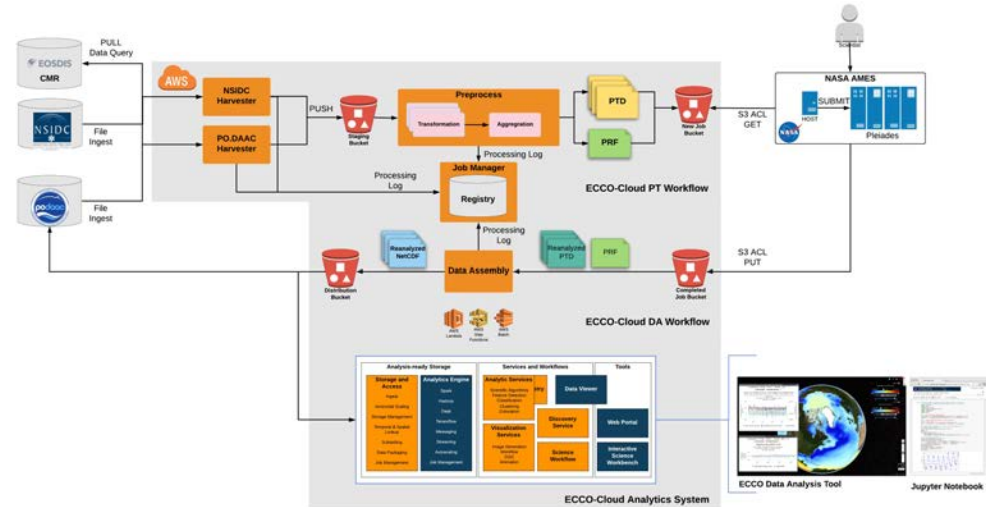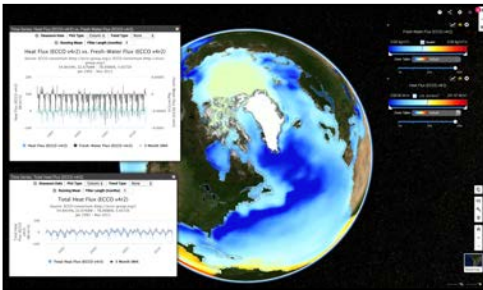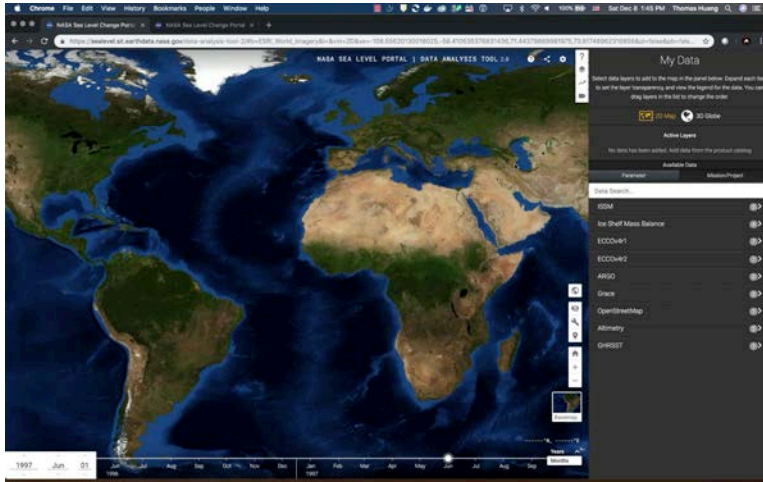- It is the first adjoint-based, multi-decadal global ocean and sea-ice state estimate
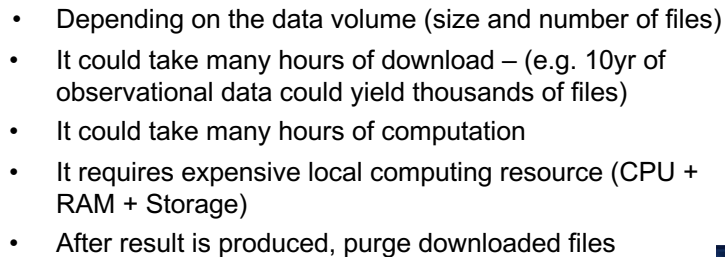


**ECCO v4r3**
- A 1-deg resolution 4D (space and time) **reconstruction** of the **1992 – 2015 global ocean** and **sea-ice state**
- With over 80 variables

- How to visualize and analyze 80 variables?
- How to compare them?
- How to compare ECCO variables with other observational variables to see how they interact?

# Estimating the Circulation and Climate of the Ocean
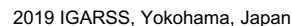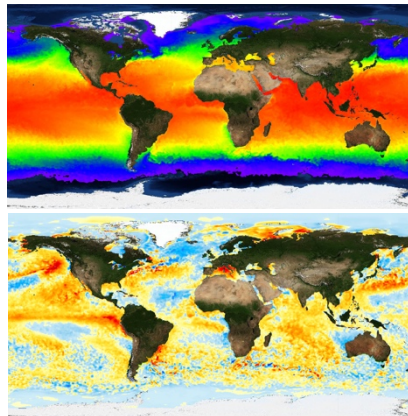## NASA ACCESS Program | PI: Patrick Heimbach; Co-Is: Ian Fenty and Thomas Huang





- The **Estimating the Circulation and Climate of the Ocean (ECCO)** global ocean state estimation system (https://ecco.jpl.nasa.gov) is the premier tool for synthesizing NASA's diverse Earth system observations into a complete physical description of Earth's time-evolving full-depth ocean and sea ice system.
- Automate generation of ECCO reanalysis products into CF-compliant netCDF products
- Integrating Amazon Cloud with NASA Ames Pleiades petascale supercomputer
- Establish ECCO Data Analysis Services and web portal for interactive visualization and analysis, and distribution
- Support multi-dimensional data visualization and analysis

# Traditional Method for Analyze Satellite Measurements

Search → Download → Compute

- Depending on the data volume (size and number of files)
- It could take many hours of download – (e.g. 10yr of observational data could yield thousands of files)
- It could take many hours of computation
- It requires expensive local computing resource (CPU + RAM + Storage)
- After result is produced, purge downloaded files

**Observation**

- Traditional methods for data analysis (time-series, distribution, climatology generation) can't scale to handle large volume, high-resolution data. They perform poorly
- Performance suffers when involve large files and/or large collection of files
- A high-performance data analysis solution must be free from file I/O bottleneck



Temporal spatial Arrays

# Processors are not Getting Faster

2004: First Pentium 4 processor with 3.0GHz clock speed

2018: Apple's MacBook Pro has clock speed of 2.7GHz

14 years later, not much has gain in raw processing power

**Modern big data architects are required to "think outside of the box". Literally!**
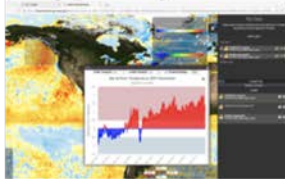
42 Years of Microprocessor Trend Data



Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California



NASA Sea Level Change Portal

Oceanographic Anomaly Detection

PO.DAAC State Of The Ocean

Hydrological Basin Analysis

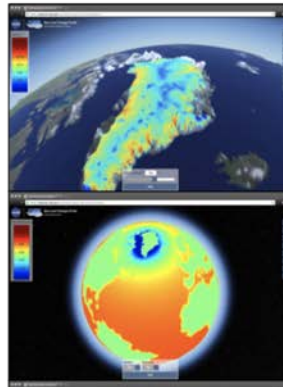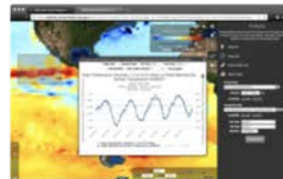Jupyter Notebook - Interactive Workbench
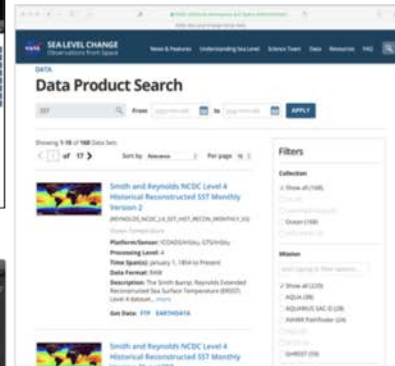
Mobile Analysis

In Situ Data Analysis
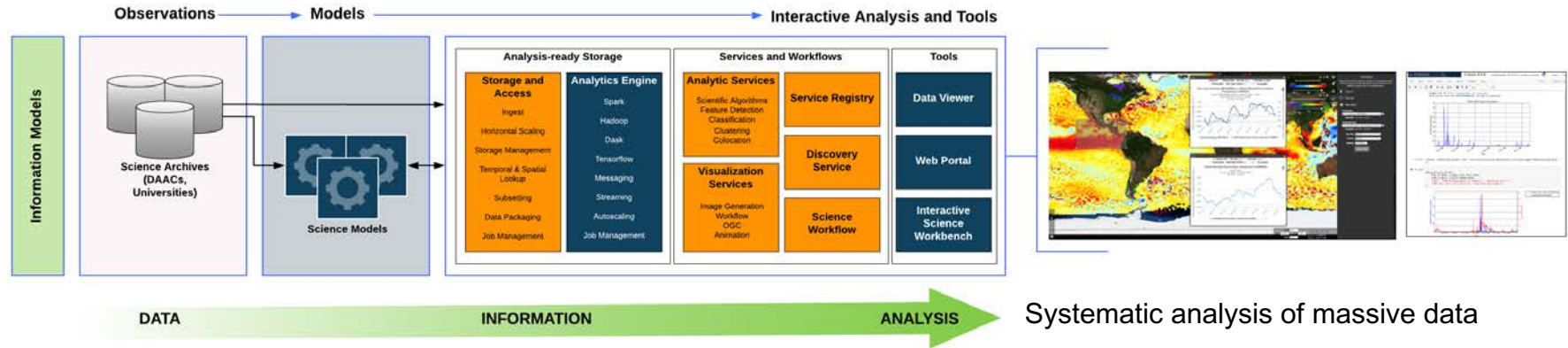
Model Simulations

Model - Observation Comparison

Integrated Search and Discovery
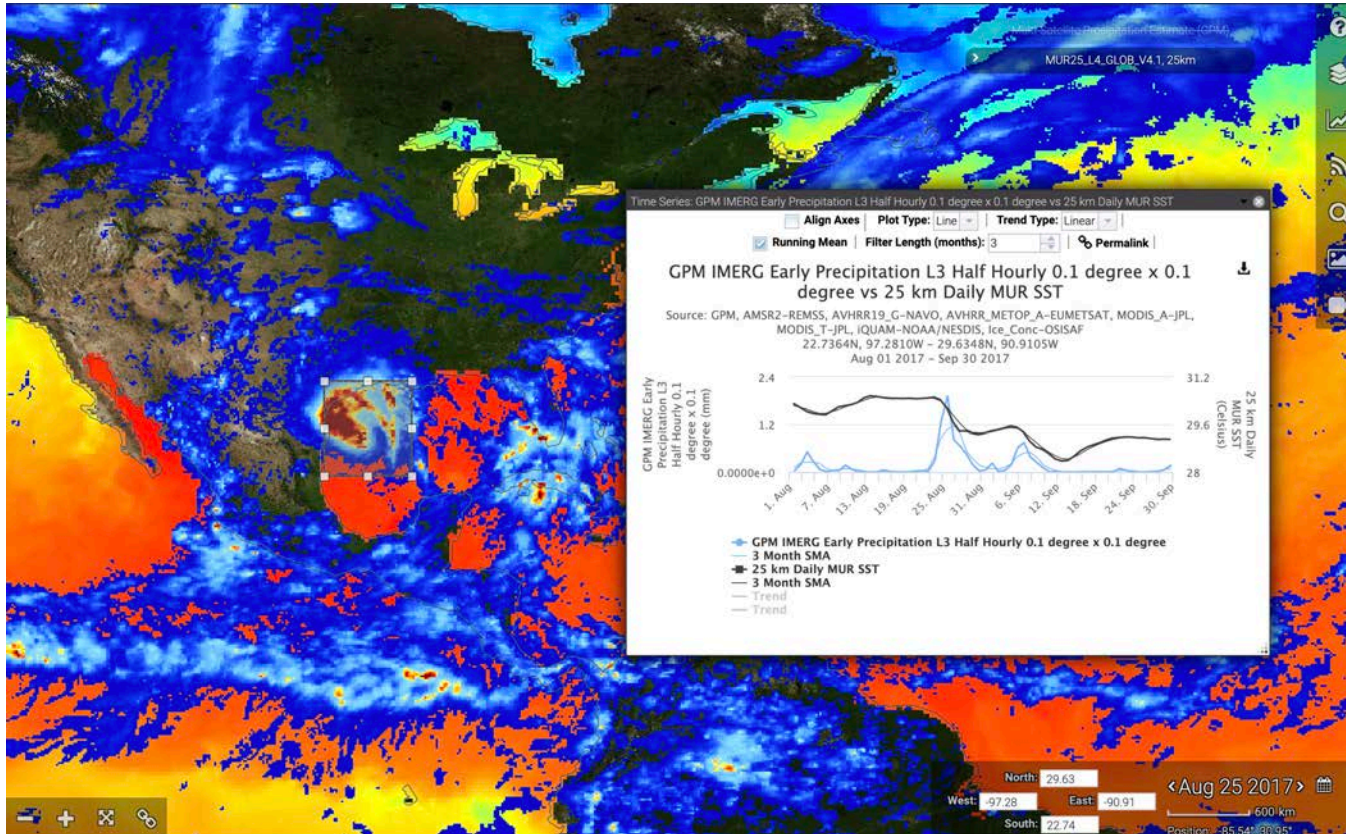
# Integrated Science Data Analytics Platform
## Creating SaaS and PaaS for Science Tools and Services



Systematic analysis of massive data

- **Integrated Science Data Analytics Platform**: an analytic center framework to provide an environment for conducting a science investigation
  - Enables the confluence of resources for that investigation
  - Tailored to the individual study area (physical ocean, sea level, etc.)
- Harmonizes data, tools and computational resources to permit the research community to focus on the investigation
- Scale computational and data infrastructures
- Shift towards integrated data analytics
- Algorithms for identifying and extracting interesting features and patterns

# Analyze Hurricane Harvey using GPM and SST
## Aug 17, 2017 – Sept. 2, 2017

# Visualize and Analyze Sea Level
## NASA Sea Level Change Portal - https://sealevel.nasa.gov
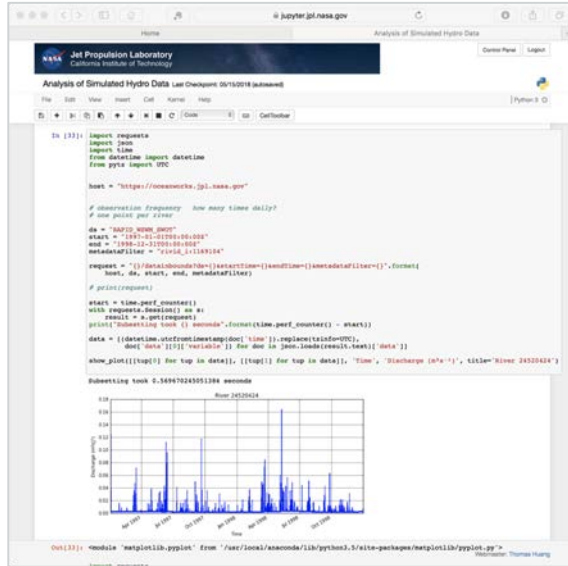


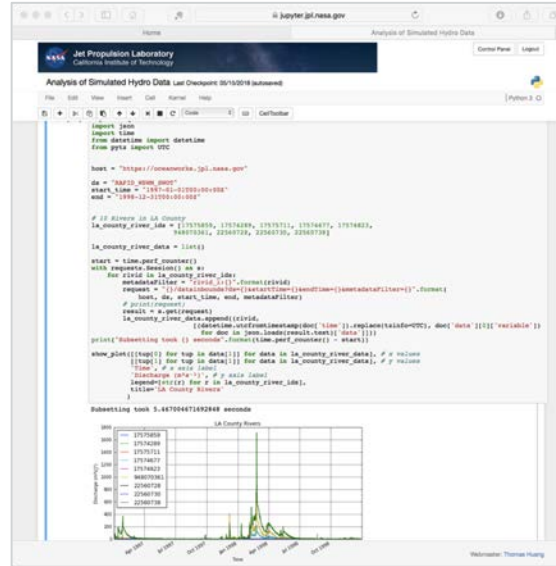Analyze *in situ* and satellite observations
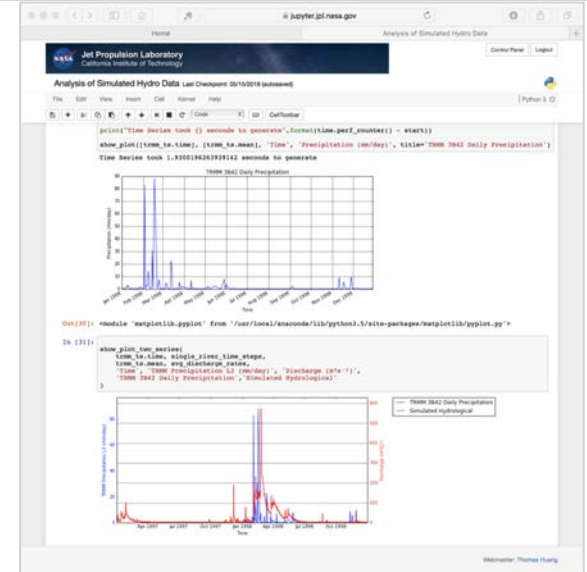


Analyze Sea Level
on mobiles

# Hydrology



Retrieval of a single river time series
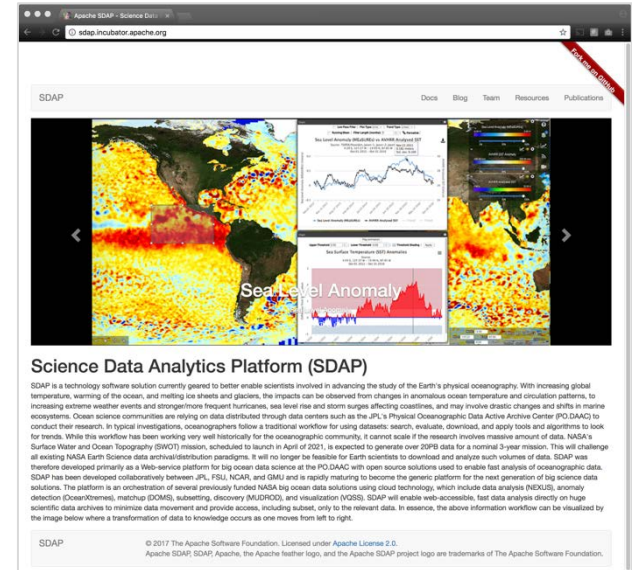


Retrieval of time series from 9 rivers



Time series coordination between TRMM and river

- Simulated hydrology data in preparation for SWOT hydrology
- **River data**: ~3.6 billion data points. 3-hour sample rate. Consists of measurements from ~600,000 rivers
- **TRMM data**: 17 years, .25deg, 1.5 billion data points
- Sub-second retrieval of river measurements
- On-the-fly computation of time series and generate coordination plot

- After more than two years of active development, on October 2017 the **NASA ESOT/AIST OceanWorks** team established Apache Software Foundation and established the **Science Data Analytics Platform (SDAP)** in the **Apache Incubator**
- Technology sharing through Free and Open Source Software (FOSS)
- Why? Further technology evolution that is restricted by projects / missions
- It is more than GitHub
  - Quarterly reporting
  - Reports are open for community review by over 6000 committers
  - SDAP has a group of appointed international mentors
- **SDAP and many of its affiliated projects are now being developed in the open**
  - Support local cluster and cloud computing platform support
  - Fully containerized using Docker and Kubernetes
  - Infrastructure orchestration using Amazon CloudFormation
  - Satellite and model data analysis: time series, correlation map,
  - In situ data analysis and collocation with satellite measurements
  - Fast data subsetting
  - Upload and execute custom parallel analytic algorithms
  - Data services integration architecture
  - OpenSearch and dynamic metadata translation
  - Mining of user interaction and data to enable discovery and recommendations



http://sdap.apache.org

APACHE INCUBATOR

# Other Applications of SDAP

National Aeronautics and Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

- **Committee of Earth Observation Satellites (CEOS) Ocean Variables Enabling Research and Applications for GEO (COVERAGE) Initiative**

- Seeks to provide **improved access** to **multi-agency ocean remote sensing data** that are **better integrated with in-situ and biological observations**, in support of **oceanographic and decision support applications** for societal benefit.

- A community-support open specification with common taxonomies, information model, and API (maybe security)

- Putting value-added services next to the data to eliminate unnecessary data movement

- Avoid data replication. Reduce unnecessary data movement and egress charges

- Public accessible RESTful analytic APIs where computation is next to the data

- Analytic engine infused and managed by the data centers perhaps on the Cloud

- Researchers can perform multi-variable analysis using any web-enabled devices without having to download files



Tools and Workbenches

Data and Information Portal

# COVERAGE Phase B



- WEkEO
  - Copernicus Data and Information Access Services (DIAS)
    1. Copernicus Data
    2. Virtual Environment and Tools
    3. User Support
  - Harmonized Data Access for Satellite data and Services
  - Virtualized infrastructure for personal sandboxes
  - Pre-configured tools
- COVERAGE Phase B
  - Establish US Node on Amazon Cloud
  - Establish EU Node on WEkEO at EUMETSAT
  - Establish COVERAGE data portal and analysis tool powered by the COVERAGE Nodes at US and EU

# PO.DAAC's SOTO

- NASA's Physical Oceanography Distributed Active Archive Center (PO.DAAC) is an element of the Earth Observing System Data and Information System (EOSDIS)
- PO.DAAC's mission is to preserve NASA's ocean and climate data and make these universally accessible and meaningful
- State of the Ocean (SOTO) is a PO.DAAC's popular visualization tool for the physical oceanography community
- SOTO v5 will be integrated with Apache SDAP and operate on the Amazon Cloud for on-the-fly data analytics



https://podaac.jpl.nasa.gov

### GRACE-FO Portal and Data Analysis Tool
GRACE/GRACE-FO Science

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

- We are developing the new GRACE Follow-On Data Science Portal – https://grace.jpl.nasa.gov

- Goals

  - Common information model

  - Unified data search and access

  - Automated, serverless data processing, analysis and image generation system

  - Integrated with Google Analytics

  - New scientific data analytics capabilities

    - Hydrological basin analysis

    - Regional – country, continent, ocean basin, etc.

    - Multivariate data analysis

  - Deploy on Amazon Web Service with auto-scaling

# Building Community-Driven Open Data and Open Source Solutions

- Deliver solutions to establish coherent platform solutions
- Embrace open source software
- Community validation
- Evolve the technology through community contributions
- Share recipes and lessons learned
- Technology demonstrations
- Host webinars, hands-on cloud analytics workshops and hackathons



2019 EGU – NASA Hyperwall





Big Data Analytics and Cloud Computing Workshop, 2017 ESIP Summer Meeting, Bloomington, IN



2019 JPL Data Science Showcase

- **The gap between visionary to pragmatists is significant**. – Geoffrey Moore
- Become an expert in the production environment and devote resources in automations
- Give project engineering team early access to the PaaS
- Deliver all technical documents and work with project system engineering
- Provide project-focused trainings



NASA's Sea Level Change Team



CEOS SIT Technical Workshop





NASA's Physical Oceanography Distributed Active Archive Center (PO.DAAC)

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

- Our new Community White Paper: "An Integrated Data Analytics Platform"
- https://www.frontiersin.org/articles/10.3389/fmars.2019.00354/full?&utm_source=Email_to_authors_&utm_medium=Email&utm_content=T1_11.5e1_author&utm_campaign=Email_publication&field=&journalName=Frontiers_in_Marine_Science&id=433796
- We are invited to discuss about "Big Ocean Science Analytics using Apache Science Data Analytics Platform" at OceanObs 2019



OCEAN OBS '19
AN OCEAN OF OPPORTUNITY
September 16-20, 2019

# ApacheCon North America

- Invited to discuss our Apache Science Data Analytics Platform (SDAP) project at the ApacheCon North America
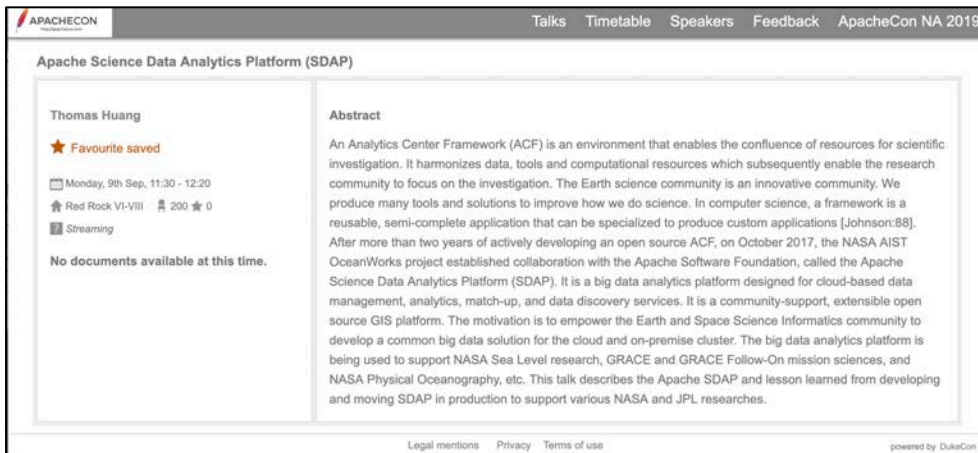
# Acknowledgement

| | | | | |
|---|---|---|---|---|
| Ed Armstrong/JPL | Maya DeBellis/JPL | Joe Jacob/JPL | David Moroni/JPL | Rob Toaz/JPL |
| Jason Barnett/LARC | Rich Doyle/JPL | Zaihua Ji/NCAR | Kevin Murphy/NASA | Vardis Tsontos/JPL |
| Andrew Bingham/JPL | Jocelyn Elya/FSU | Yongyao Jiang/GMU | Charles Norton/JPL | Suresh Vannan/JPL |
| Carmen Boening/JPL | Ian Fenty/JPL | Felix Landerer/JPL | Jean-Francois Piolle/IFREMER | Jorge Vazquez/JPL |
| Mark Bourassa/FSU | Kevin Gill/JPL | Yun Li/GMU | Nga Quach/JPL | Ou Wang/JPL |
| Mike Chin/JPL | Frank Greguska/JPL | Eric Lindstrom/NASA | Brandi Quam/NASA | Brian Wilson/JPL |
| Marge Cole/NASA | Patrick Heimbach/UT Austin | Mike Little/NASA | Shawn Smith/FSU | Steve Worley/NCAR |
| Tom Cram/NCAR | Ben Holt/JPL | Chris Lynnes/NASA | Ben Smith/JPL | Elizabeth Yam/JPL |
| Dan Crichton/JPL | Thomas Huang/JPL | Lewis McGibbney/JPL | Adam Stallard/FSU | Phil Yang/GMU |

# In Summary

- **You've got to think about big things while you're doing small things, so that all the small things go in the right direction** – Alvin Toffler
- Focus on end-to-end data and computation architecture, and the total cost of ownership
- JPL Strategy is to drive Data Science into the fabric of JPL by
    - Launching cross-institution pilots
    - Building a trained workforce
    - Linking to the mission-science data lifecycle
- Invest in Interactive Analytics that simplifies the integration of *multiple* Earth observing remote sensing instruments; comparison against models

National Aeronautics and Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California



IGARSS 2019

Save the Date!
APACHECON
North America
20 YEARS OF APACHE
Las Vegas, Nevada
September 9-12, 2019

OCEAN OBS'19
An Ocean of Opportunity

Presentation | Evaluation | Collaboration

**Thomas Huang**

thomas.huang@jpl.nasa.gov

Jet Propulsion Laboratory

California Institute of Technology